

日期: 12 / 17

## Mini-batch 梯度下降.

$m = 5000, 000.$

5000  $\uparrow$  batch, batch-size = 1000.

$x^{(t)}$ ,  $y^{(t)}$ .

(1000  $\uparrow$ ) (1000  $\uparrow$ ).

for  $t=1$  to 5000:

对  $x^{(t)}$ ,  $y^{(t)}$  做梯度下降.

$$\text{cost} = J^{(t)} = \frac{1}{1000} \cdot \sum_{i=1}^{1000} L(y^{(i)}, y^{(i)}) + \frac{\lambda}{2 \times 1000} \|w^{(t)}\|_F^2$$

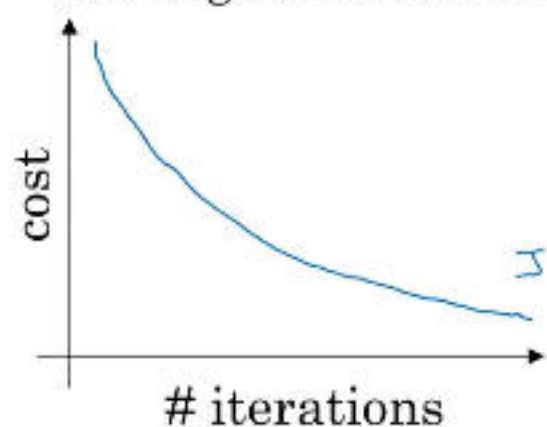
Back Propagate.

Gradient Descent.

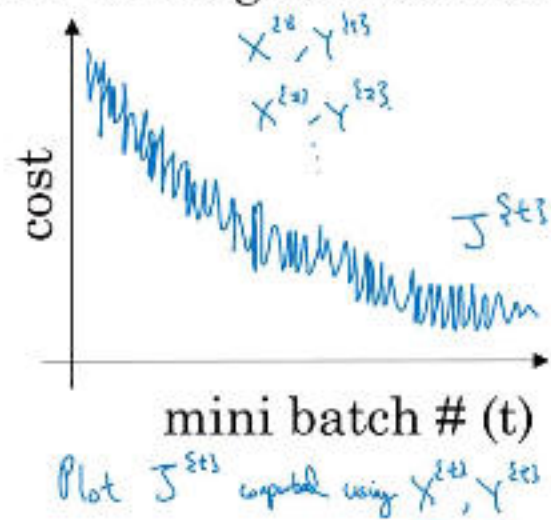
一个 epoch: 遍历  $\{x, y\}$  一次, 已经做了 5000 次梯度下降.  
但每一次都较快.

### Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent



Andrew Ng

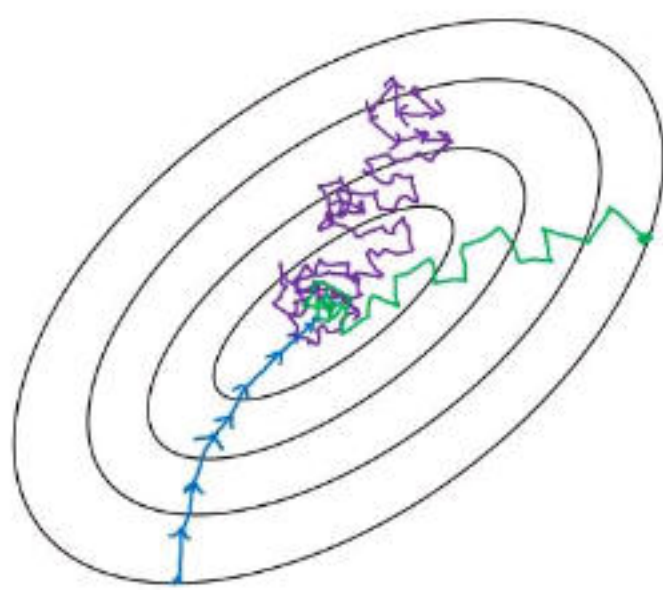
日期: 12 / 17

## Choosing your mini-batch size

→ If mini-batch size =  $m$  : Batch gradient descent.  $(X^{(1)}, Y^{(1)}) = (X, Y)$ .

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is its own mini-batch.  
 $(X^{(1)}, Y^{(1)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$  mini-batch.

In practice: Somewhere in-between 1 and  $m$



Stochastic  
gradient  
descent

Loss spikes  
due to vectorization

In-between  
(mini-batch size  
not too big/small)

Faster learning.

- Vectorization.  
(~1000)
- Make passes without  
processing entire training set.

Batch  
gradient descent  
(mini-batch size =  $m$ )

Too long  
per iteration

Andrew Ng

(Gradient  
Descent).

batch size = 1 : 一个样本做一次梯度下降,  $\Rightarrow$  随机 GD.

batch size =  $m$  : 所有样本, 全局优化, 可能很慢 ( $m$  大)

batch size 常取  $2^6, \dots, 2^{10}$ . 2 进制, 但方向性强.



日期: 12/17

## 指数加权平均.

### Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$\begin{aligned} v_{100} &= 0.9v_{99} + 0.1\theta_{100} \\ v_{99} &= 0.9v_{98} + 0.1\theta_{99} \\ v_{98} &= 0.9v_{97} + 0.1\theta_{98} \\ &\dots \end{aligned}$$

$$\begin{aligned} v_{100} &= 0.1\theta_{100} + 0.9(0.1\theta_{99} + 0.9v_{98}) \\ &= 0.1\theta_{100} + 0.1 \times 0.9 \theta_{99} + 0.1(0.9)^2 \theta_{98} + 0.1(0.9)^3 \theta_{97} + \dots \\ 0.9^{50} &\approx 0.35 \approx \frac{1}{e} \end{aligned}$$

展开:  $0.1 \uparrow$

$\sum = 1 - \beta$

$0.1\theta_{99} + 0.9v_{98}$

$(1-\epsilon)^{\frac{1}{\epsilon}} = \frac{1}{e}$

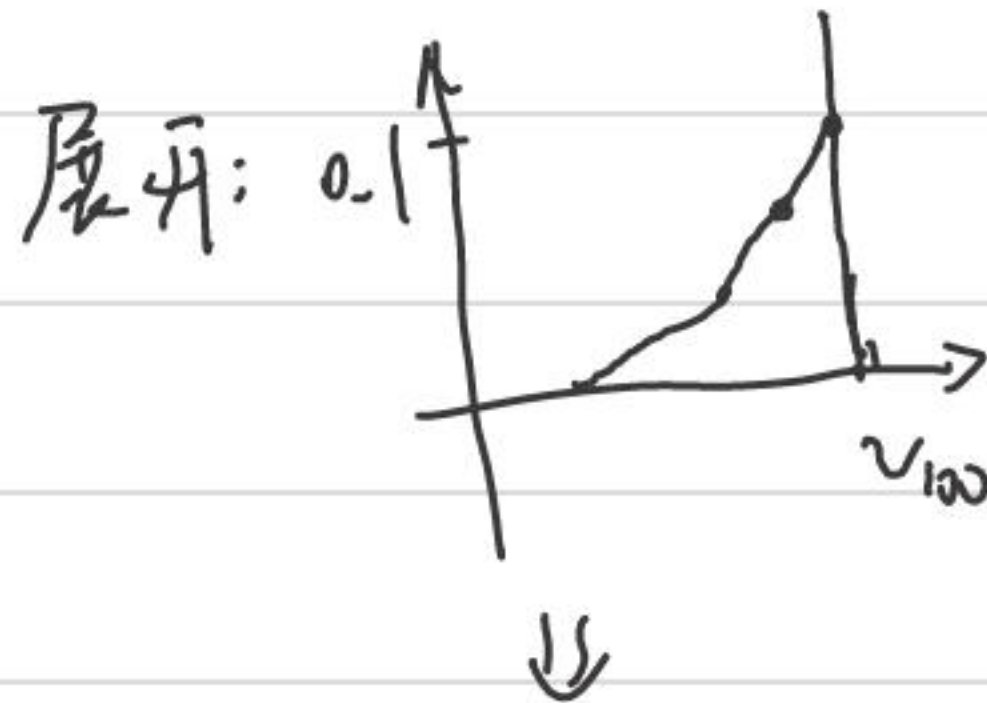
$\epsilon = 0.02 \rightarrow 0.98^{50} \approx \frac{1}{e}$

Andrew Ng

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_0 = \theta_0$$



$\therefore$  近似平均了  $\frac{1}{1-\beta}$  天的数据.

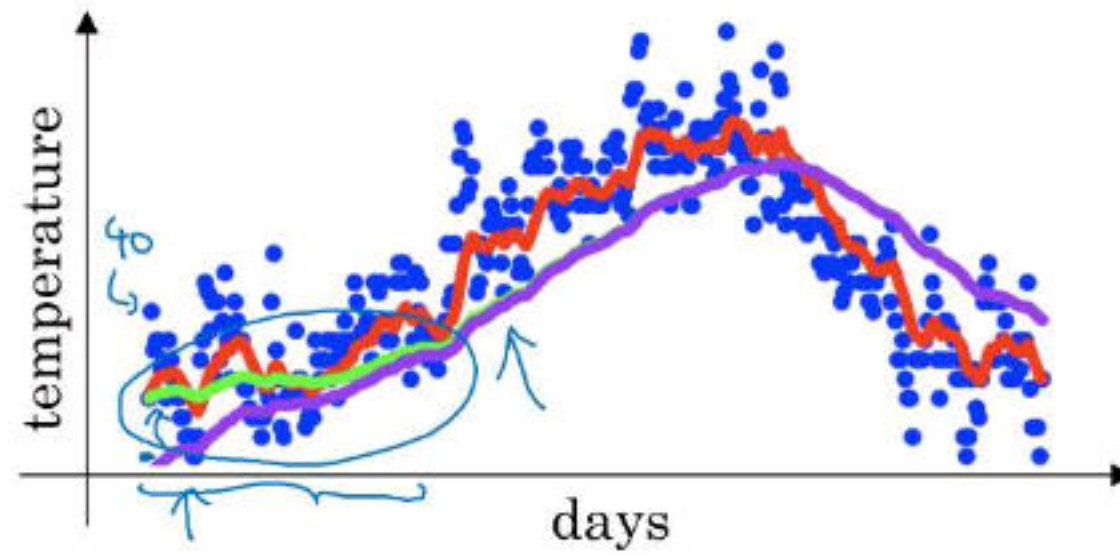
指数 Decay.

代码只需维护变量  $v$ .

问题:

日期: 12 / 17

## Bias correction



$$\rightarrow v_t = \beta v_{t-1} + (1 - \beta)\theta_t$$

$$v_0 = 0$$

$$v_1 = \cancel{0.98 v_0} + 0.02 \theta_1$$

$$\begin{aligned} v_2 &= 0.98 v_1 + 0.02 \theta_2 \\ &= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2 \\ &= 0.0196 \theta_1 + 0.02 \theta_2 \end{aligned}$$

$$\frac{v_t}{1 - \beta^t}$$

$$t=2: 1 - \beta^t = 1 - (0.98)^2 = 0.0396$$

$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$

Andrew Ng

$$v_0 = 0$$

问题:

$$v_0 = 0$$

$$v_1 = 0.98 v_0 + 0.02 \theta_1 = \frac{0.02 \times 40}{\delta}$$

$$v_2 = 0.98 v_1 + 0.02 \theta_2 = \frac{0.98 \times 0.02 \theta_1 + 0.02 \theta_2}{\delta} \quad \text{Bias.}$$

$$\frac{v_t}{1 - \beta^t}$$

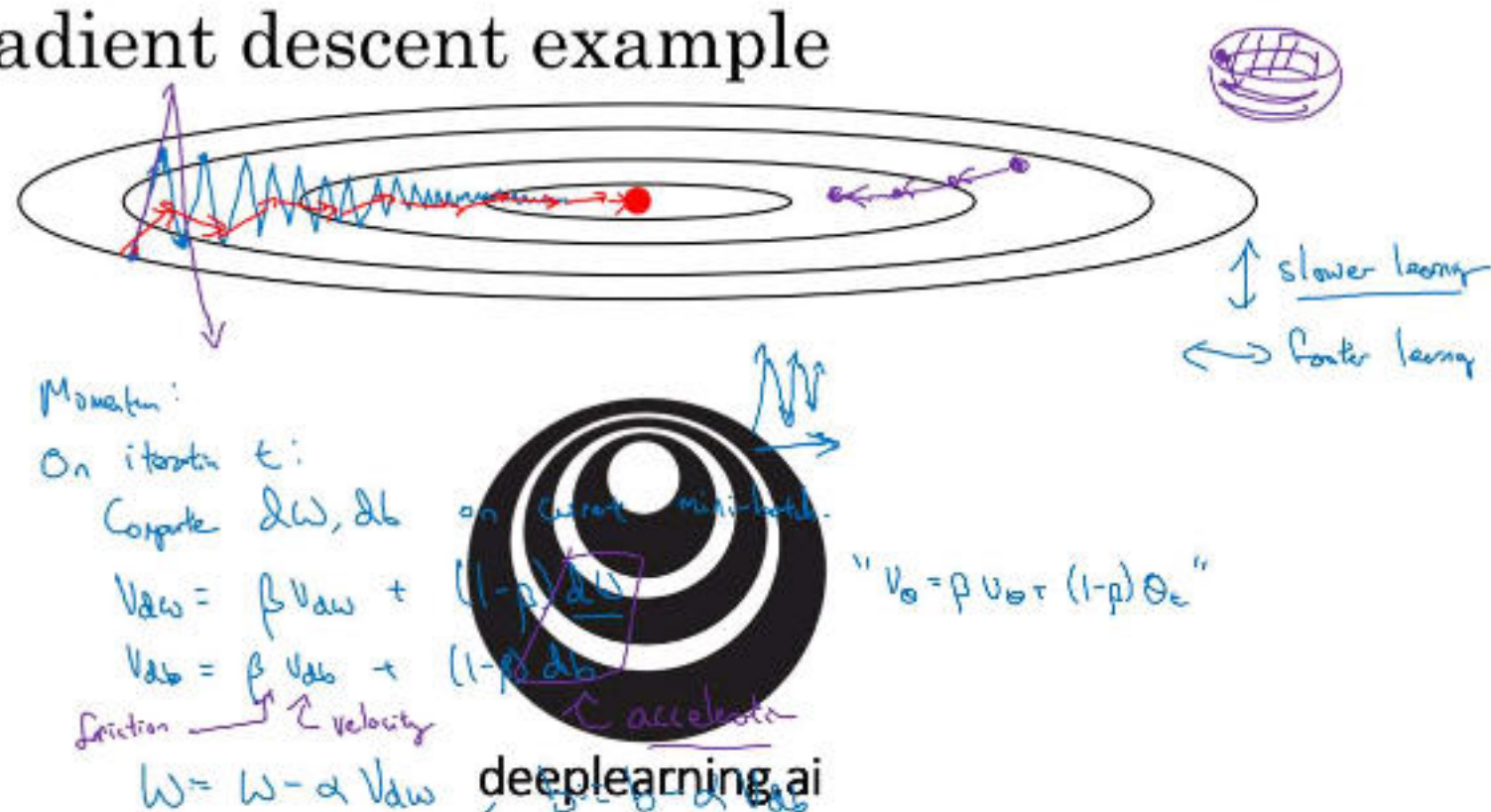
$$\beta \rightarrow 1, \quad \beta^t \rightarrow 0 \quad (t \uparrow)$$



日期: 12 / 17

# 三种优化算法:

## Gradient descent example



Andrew Ng

a): Momentum:

$$v_{dw} = \beta v_{dw} + (1-\beta) dw$$

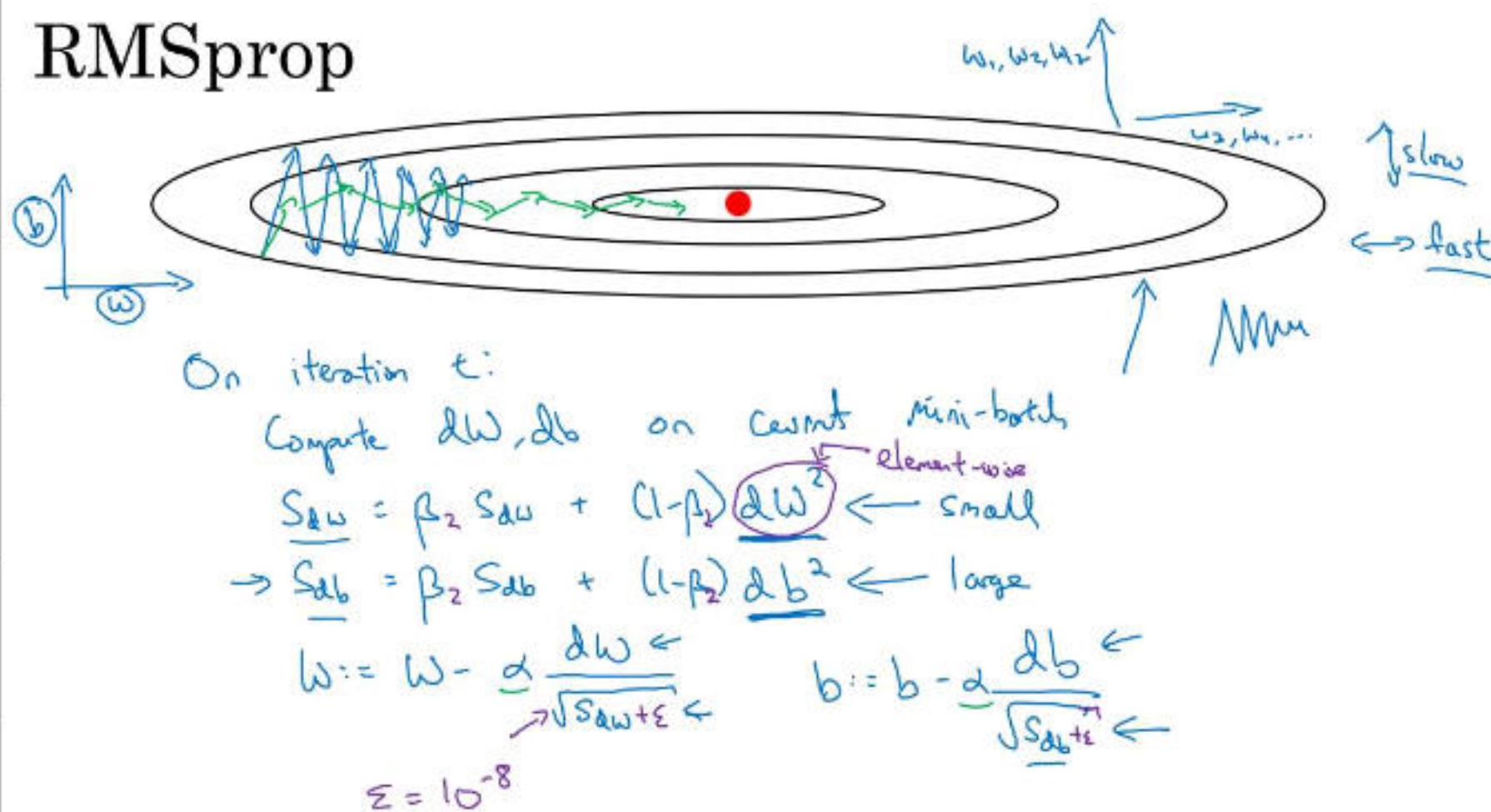
对  $dw, db$  做

$$v_{db} = \beta v_{db} + (1-\beta) db$$

指数加权平均。

垂直方向抵消, 减少无谓的摆动。

## RMSprop



Andrew Ng



b): Root Mean Square prop.

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dw^2$$

element-wise.

$dw \uparrow, S_{dw} \uparrow$

$$S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

$S_{db} \downarrow$

$$W = W - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}$$

$$b = b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$



日期: 12 / 17

## Adam optimization algorithm

$V_{dw}=0, S_{dw}=0, V_{db}=0, S_{db}=0$

On iteration  $t$ :

Compute  $dw, db$  using current mini-batch

$V_{dw} = \beta_1 V_{dw} + (1-\beta_1)dw, V_{db} = \beta_1 V_{db} + (1-\beta_1)db \leftarrow \text{"momentum"} \beta_1$

$S_{dw} = \beta_2 S_{dw} + (1-\beta_2)dw^2, S_{db} = \beta_2 S_{db} + (1-\beta_2)db^2 \leftarrow \text{"RMSprop"} \beta_2$

$yhat = np.array([.9, 0.2, 0.1, .4, .9])$

$V_{dw}^{corrected} = V_{dw} / (1-\beta_1^t), V_{db}^{corrected} = V_{db} / (1-\beta_1^t)$

$S_{dw}^{corrected} = S_{dw} / (1-\beta_2^t), S_{db}^{corrected} = S_{db} / (1-\beta_2^t)$

$W := W - \alpha \frac{V_{dw}^{corrected}}{\sqrt{S_{dw}^{corrected} + \epsilon}}, b := b - \alpha \frac{V_{db}^{corrected}}{\sqrt{S_{db}^{corrected} + \epsilon}}$

Andrew Ng

c): Adaptive momentum.

对  $dw, db$  指数加权.

同时对修正后的  $V_{dw}, V_{db}$  做 "RMS".

两种优化方法的叠加.